

# Dealing with Microarray Data

NIDDK Meeting  
10 July 2001



## The Golden Age of Genomics

- ~50 Microbial Genomes have been sequenced, at least 50 more are on the way
- Yeast, *C. elegans*, *Arabidopsis*, *Drosophila* and other Eukaryotic models are finished or well advanced
- A “working draft” of the Human Genome Sequence is now available, with mouse and rat to follow
- More than 8,000,000 ESTs are available; more than 3,600,000 from humans

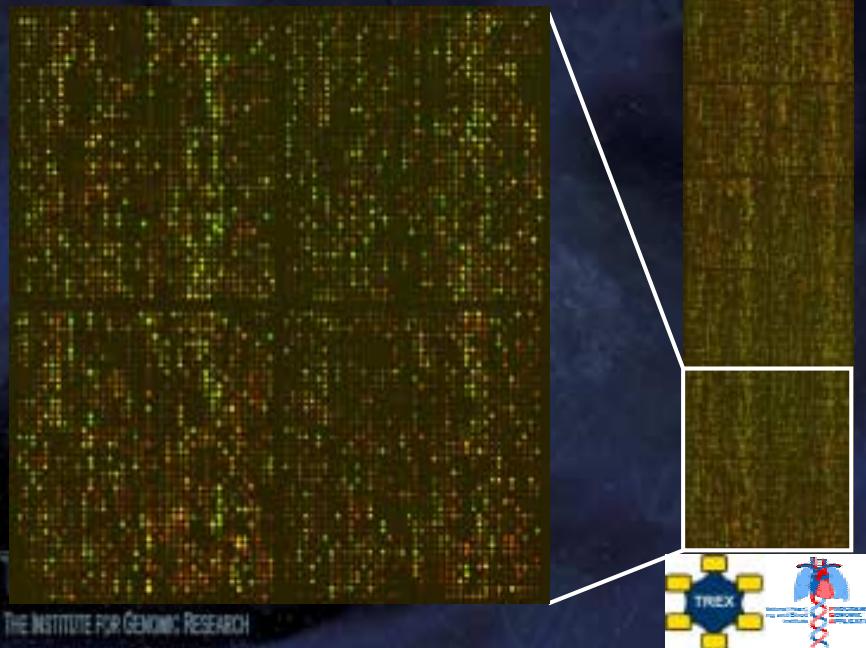


Science is built with facts as a house is with stones – but a collection of facts is no more a science than a heap of stones is a house.

– Jules Henry Poincare



## Hybridization to a 19,200 Element Human Array



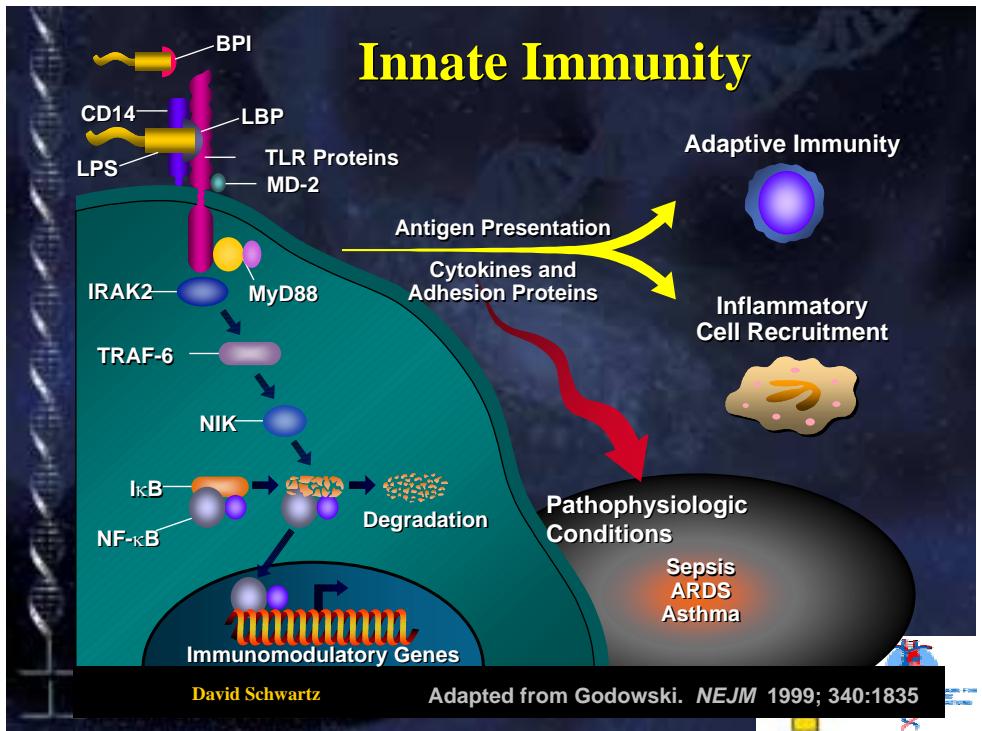
## **Named Differentially Expressed Genes**

# Microarray Expression Profiling of Rodent Models of Human Disease

## **The Theme of our PGA:**

# *Examining Gene/Environment Interactions in Rodent Models of Human Disease Using cDNA Microarrays to link Phenotype to Genotype*

See <<http://pga.tigr.org>> for information, data, tools, and Visiting Scientist program



## Challenges in HLB Research

- Identifying appropriate models
  - Identifying genes involved in HLB disorders
  - Placing these genes into relevant pathways
  - Separating genetic and environmental components
  - Developing an understanding of the mechanisms

## Expression Profiling: cDNA Microarrays

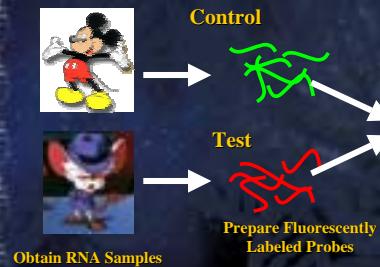
- **The Promise:**
  - Genome Scale Data on Gene Expression Patterns
- **The Challenges:**
  - Deduce Gene Function from Expression
  - Elucidate Functional Pathways from Function and Expression
- **The Solutions:**
  - Development of Robust Protocols for large-scale correlation studies
  - Relational Database for Data Storage
  - Data Visualization and Mining Tools
  - Tools for modeling pathways



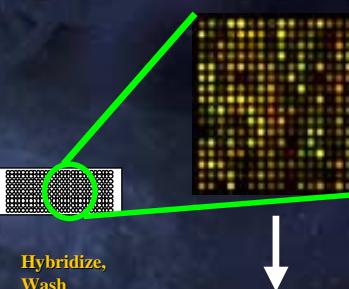
THE INSTITUTE FOR GENOMIC RESEARCH



## Microarray Overview II



Measure  
Fluorescence  
in 2 channels  
red/green



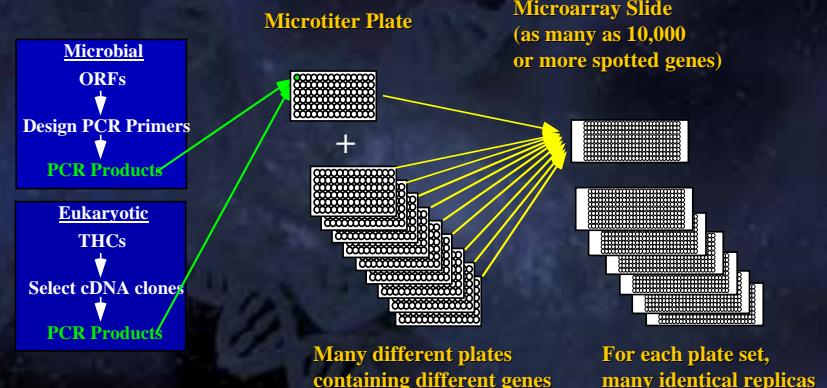
Analyze the data  
to identify  
differentially  
expressed genes



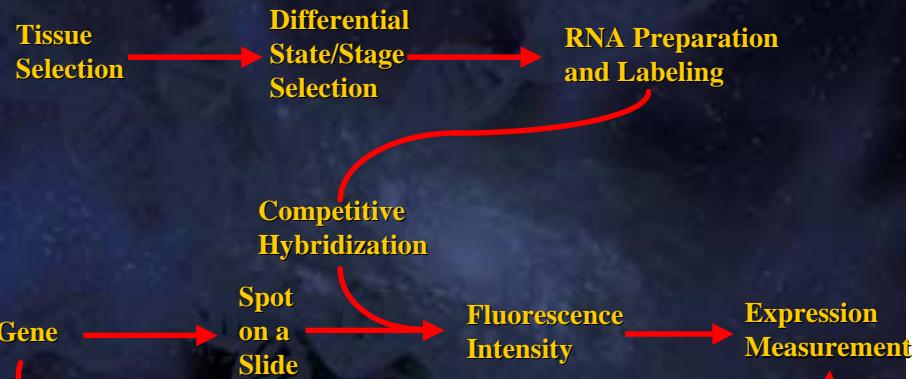
THE INSTITUTE FOR GENOMIC RESEARCH



## Microarray Overview I



## Microarray Expression Analysis



## The Beast: Microarray Robot from Intelligent Automation

<http://www.ias.com>

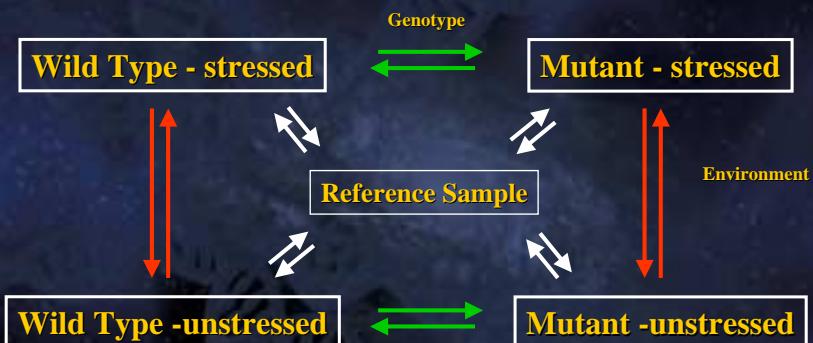


TIGR

THE INSTITUTE FOR GENOMIC RESEARCH



## Basic Experimental Paradigm



TIGR

THE INSTITUTE FOR GENOMIC RESEARCH

## Axon GenePix 4000B Scanner

[http://www.axon.com/GN\\_GenePix4000.html](http://www.axon.com/GN_GenePix4000.html)



## Developmental Goals

- Tools for linking Genes, ESTs, cDNAs, and Sequences
- Microarray Quality Control Protocols and Reagents
- Novel Analysis Techniques and Tools
- Phenotyping Pipelines
- Microarray Resources in Mouse and Rat
- Expression Profiles for Disease Phenotypes of HLBS interest

TIGR

THE INSTITUTE FOR GENOMIC RESEARCH



# Developmental Goals

- Tools for linking Genes, ESTs, cDNAs, and Sequences
- Microarray Quality Control Protocols and Reagents
- Novel Analysis Techniques and Tools
- Phenotyping Pipelines
- Microarray Resources in Mouse and Rat
- Expression Profiles for Disease Phenotypes of HLBS interest

TIGR

THE INSTITUTE FOR GENOMIC RESEARCH



## TIGR Gene Indices home page

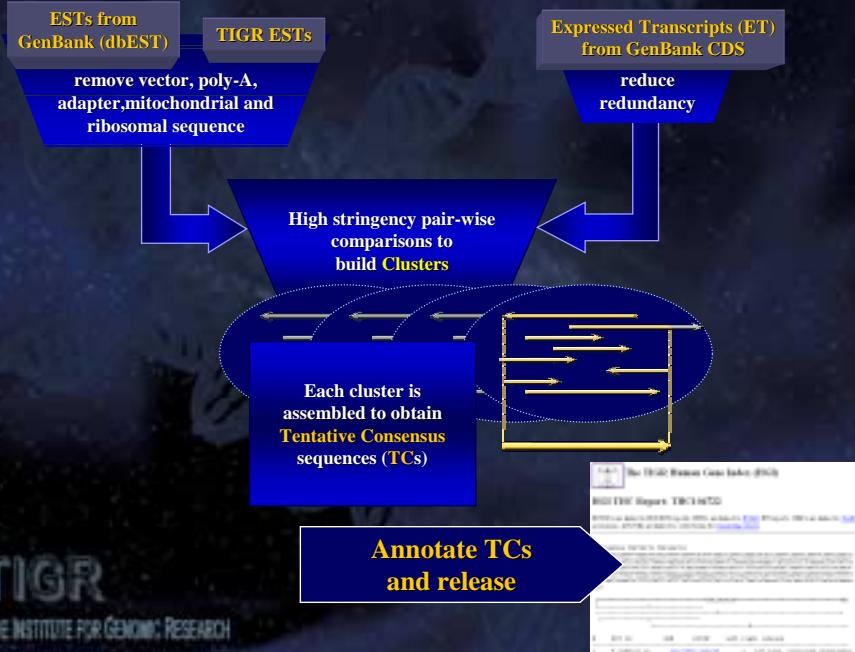
[www.tigr.org/tdb/tgi.shtml](http://www.tigr.org/tdb/tgi.shtml)



TIGR

THE INSTITUTE FOR GENOMIC RESEARCH

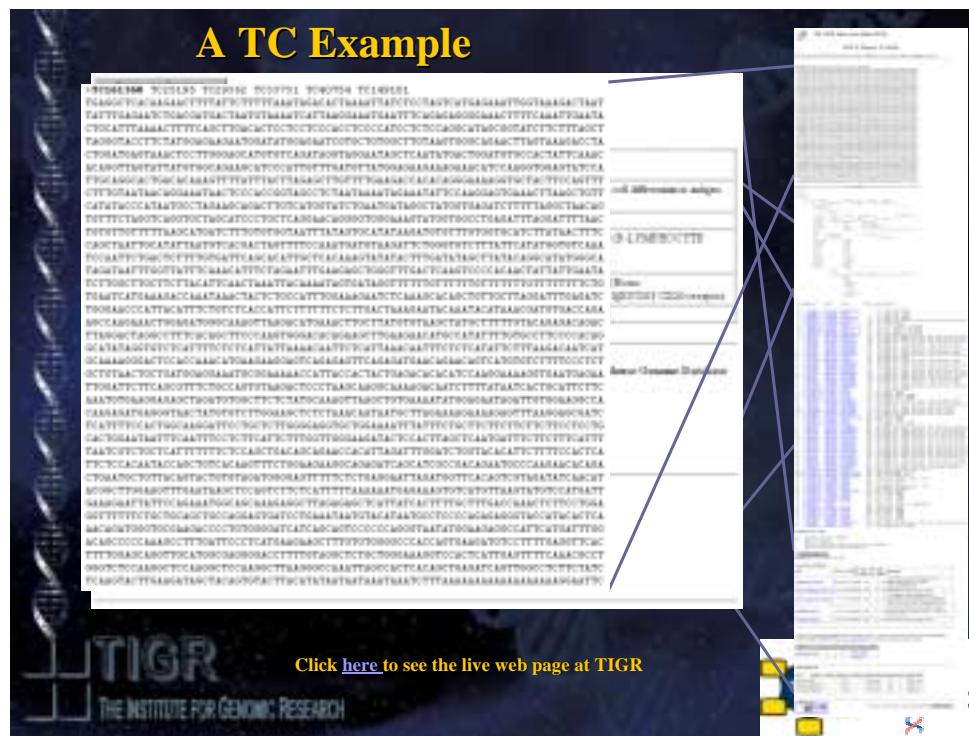
## Gene Index Assembly process



## The Mouse Gene Index <<http://www.tigr.org/tdb/mgi>>



## A TC Example



## Ancillary Information

- Annotation based on gene content
- Annotation based on protein similarity searches
- “Opposite End” information
- Links to mapping information
- Links to Tentative Orthologue Groups (TOGs)
  
- Links to Gene Ontology (GO) Terms (in release)
- Links to EC Numbers (in release)
- Links to completed genomic sequence (in release)
  
- Links to species-specific databases – MGD, RGD (in testing)



### Can ESTs be used to identify orthologues?

Property	Rat-Human	Mouse-Human	Mouse-Rat
	Mean (SD)	Mean (SD)	Mean (SD)
<b>5'UTR</b>			
% Identity	68.4 (13.0)	69.7 (12.9)	84.5 (12.9)
Mutation distance K	0.486 (0.260)	0.493 (0.273)	0.212 (0.224)
<b>CDS</b>			
% Identity (protein)	88.0 (11.8)	86.4 (12.3)	94.5 (6.3)
% Identity (DNA)	85.9 (6.0)	85.2 (6.5)	93.8 (3.2)
Syn. distance K <sub>s</sub>	4.60 (0.245)	0.468 (0.169)	0.166 (0.061)
Nonsyn. distance K <sub>a</sub>	0.078 (0.095)	0.090 (0.102)	0.031 (0.040)
<b>3' UTR</b>			
% Identity	70.1 (11.4)	71.0 (12.2)	86.3 (8.9)
Mutation distance K	0.435 (0.212)	0.447 (0.225)	0.164 (0.152)

Analysis of 1,880 rodent-human orthologue pairs:

1,212 rat-human pairs, 1,138 mouse-human pairs, and 470 shared by all three.  
Makalowski and Boguski (1998) *Proc. Natl. Acad. Sci. USA* 95, 9407-9412

### Can ESTs be used to identify orthologues?

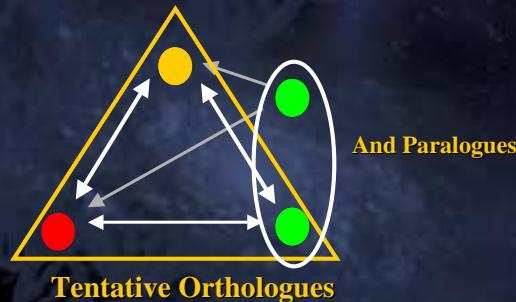
Property	Rat-Human	Mouse-Human	Mouse-Rat
	Mean (SD)	Mean (SD)	Mean (SD)
<b>5'UTR</b>			
% Identity	68.4 (13.0)	69.7 (12.9)	84.5 (12.9)
Mutation distance K	0.486 (0.260)	0.493 (0.273)	0.212 (0.224)
<b>CDS</b>			
% Identity (protein)	88.0 (11.8)	86.4 (12.3)	94.5 (6.3)
% Identity (DNA)	85.9 (6.0)	85.2 (6.5)	93.8 (3.2)
Syn. distance K <sub>s</sub>	4.60 (0.245)	0.468 (0.169)	0.166 (0.061)
Nonsyn. distance K <sub>a</sub>	0.078 (0.095)	0.090 (0.102)	0.031 (0.040)
<b>3' UTR</b>			
% Identity	70.1 (11.4)	71.0 (12.2)	86.3 (8.9)
Mutation distance K	0.435 (0.212)	0.447 (0.225)	0.164 (0.152)

Analysis of 1,880 rodent-human orthologue pairs:

1,212 rat-human pairs, 1,138 mouse-human pairs, and 470 shared by all three.  
Makalowski and Boguski (1998) *Proc. Natl. Acad. Sci. USA* 95, 9407-9412



## Building TOGs: Reflexive, Transitive Closure



**TIGR**  
THE INSTITUTE FOR GENOMIC RESEARCH



## The TIGR Gene Indices <<http://www.tigr.org/tdb/tdb/tgi.html>>



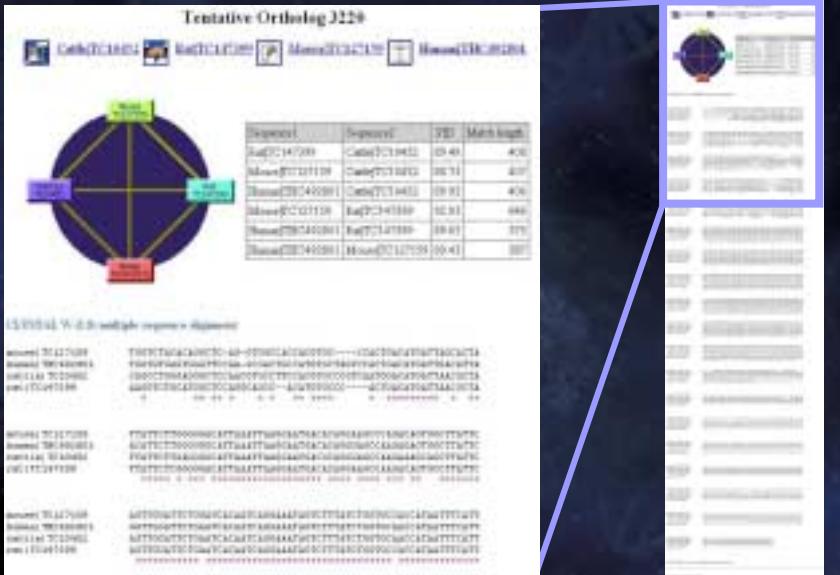
TIGR Orthologous Gene Alignments



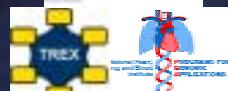
The TIGR Orthologous Gene Alignments (TOGA) database provides links between candidate orthologs identified using the Tentative Consensus (TC) sequences that comprise the TIGR Gene Indexes. TOGA currently provides links between the Human, Mouse, and Rat Gene Indexes. Future releases will incorporate additional species.



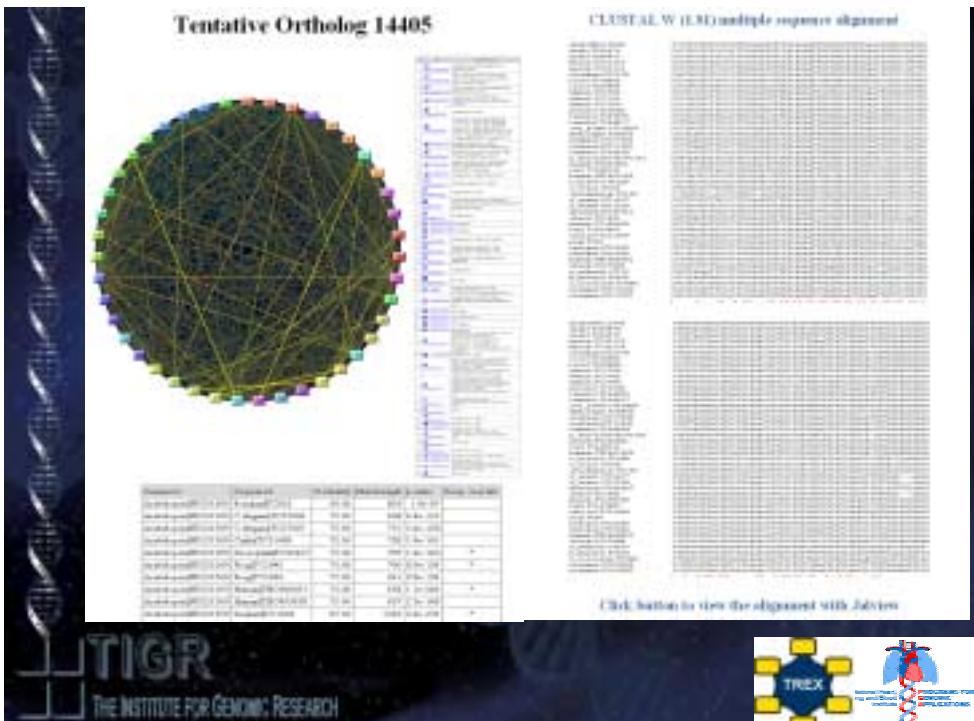
## TOGA: An Sample Alignment: bithoraxoid-like protein



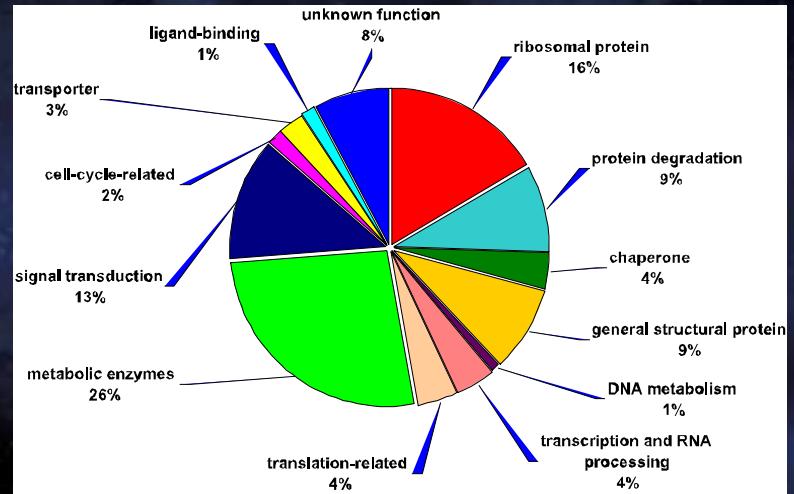
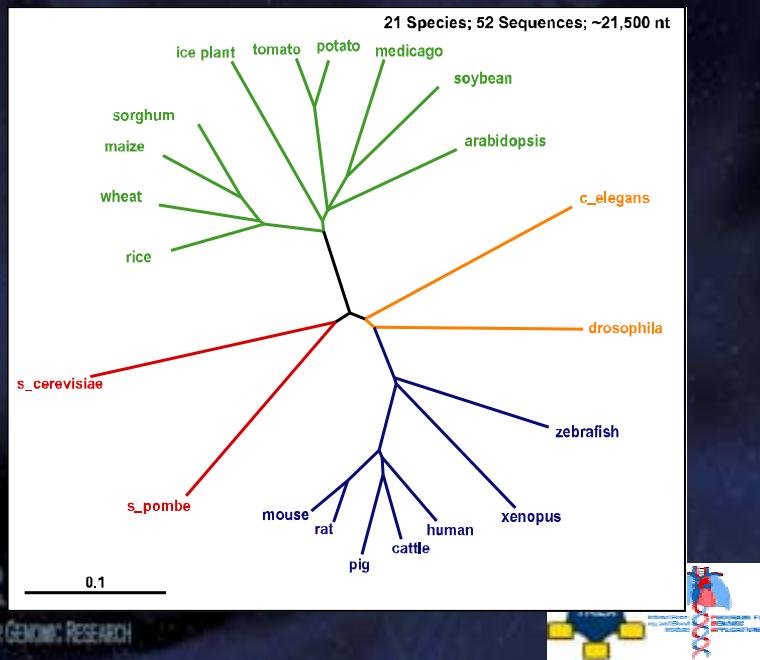
**TIGR**  
THE INSTITUTE FOR GENOMIC RESEARCH



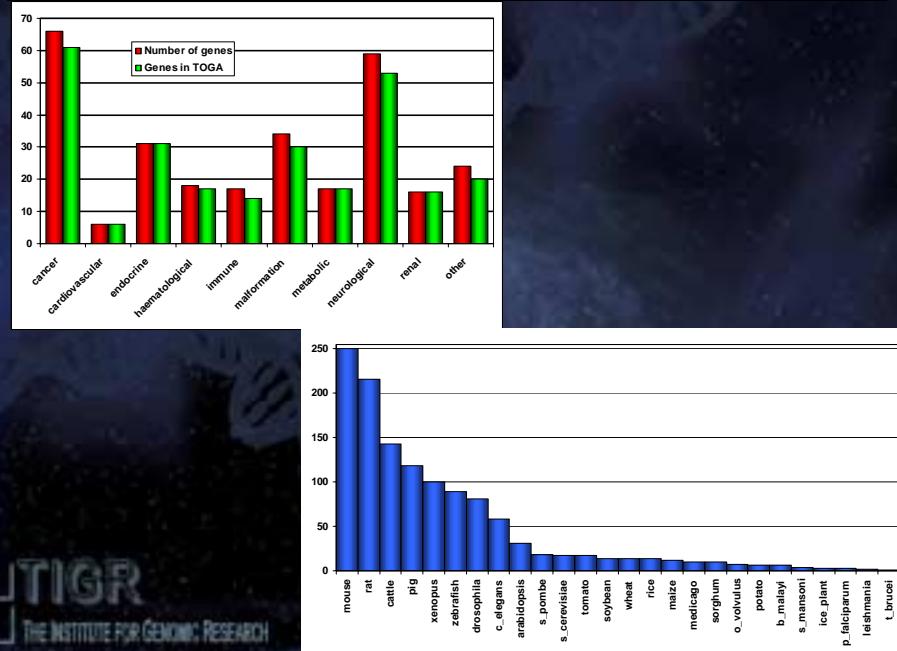
## Tentative Ortholog 14405



## Relationships in TOGA



## Representation of 288 Human Disease Genes in TOGA



## *Arabidopsis thaliana*: Alignments with Chromosome II



<http://www.tigr.org/tdb/at/alignTC.html>

## Links from Public Clone Sets to TGI

**NHLBI Tools**

Our PGA is committed to making the tools, reagents, and protocols developed through this project to the research community.

- Links From Common Microarray Resources to the TIGR Gene Indices
  - **CGPSA**
    - **Mouse**: [SMAF](#), [RNA](#)
    - **Rat**
    - **Arabidopsis**
    - **Drosophila**
      - **Mouse**: [Mouse Genome Chip Site](#)
      - **Human**: [Human Genome Chip Site](#)
  - **Microarray protocols**, including protocols and links to microarray analysis software developed at TIGR.
  - Access to the microarray generated by our TGA will be provided through a [Web interface](#).
  - Access to our expertise will be provided through a [series of online consultation sessions](#).
  - Documentation for our TGA will made available within thirty days of generation.
  - Additional resources that may be of interest include:
    - The TIGR Gene Indices
    - The TIGR Mouse Gene Aligned Tissue (TIGRA) database
    - TIGRdbm
    - The Mouse Genome Database
    - The Rat Genome Database

What's New | Discourse | Help | Terms of Use | Contact Us | Links | Home

<http://pga.tigr.org/tools.shtml>




## Links from Public Clone Sets to TGI

**NIA Mouse 15k cDNA collection**

The library consists of clones for over 10,000 mouse genes ([GenBank ID: M32000-M32999](#)). These have been measured from more than 1,000,000 individual mice. T1 RNA-fingerprinting was used for clone assignment, and hierarchical clustering

**Description**

There are 15,000 clones in this collection. [Download](#) [View](#)

ProbeID	Library ID	ProbeID	Tissue ID	15KID	TIGR Expression	15KID	TIGR Expression	15KID	TIGR Expression	15KID	TIGR Expression
CP00129-0	B0000000	AA000000	Mm-1129		NA	AA000000	NA	AA000000	NA	AA000000	NA
CP00130-0	B0000000	AA000000	Mm-1130		NA	AA000000	NA	AA000000	NA	AA000000	NA
CP00131-0	B0000000	AA000000	Mm-2129		NA	AA000000	NA	AA000000	NA	AA000000	NA
CP00132-0	B0000000	AA000000	Mm-2130		NA	AA000000	NA	AA000000	NA	AA000000	NA
CP00133-0	B0000000	AA000000	Mm-2131		NA	AA000000	NA	AA000000	NA	AA000000	NA

<http://pga.tigr.org/tools.shtml>

**TIGR**  
THE INSTITUTE FOR GENOMIC RESEARCH

## RESOURCERER

**RESOURCERER**

RESOURCERER provides a platform for comparing two TGA data sets. It can compare two TGA data sets and [GeneChips](#). RESOURCERER also allows comprehensive feature matching and enables users to query gene expression using the [TIGRdbm](#) database.

Take a single measure or an "Orphan" (i.e. "Unknown") TGA dataset for comparison.

To compare two measures, select both, choose the basic comparison (T15K or T100K), and the type of comparison (plotted).

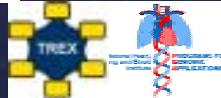
[Common gene selection](#)

**Data Set A:**  [Select](#)

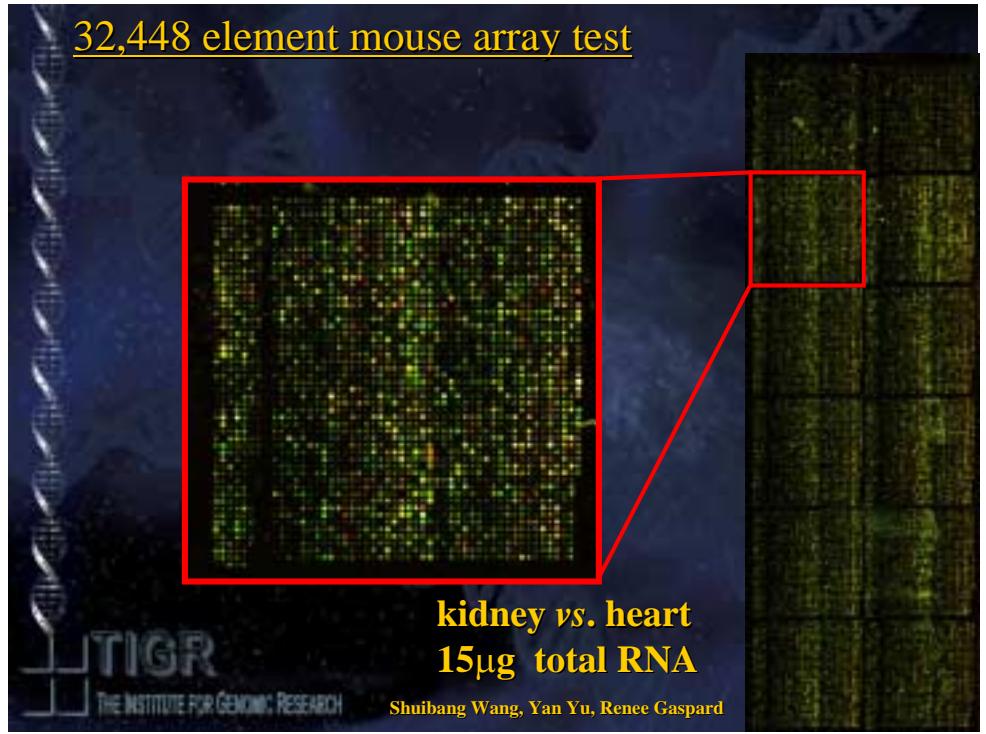
**Data Set B:**  [Select](#)  [T15K](#)  [T100K](#)

What's New | Discourse | Help | Terms of Use | Contact Us | Links | Home

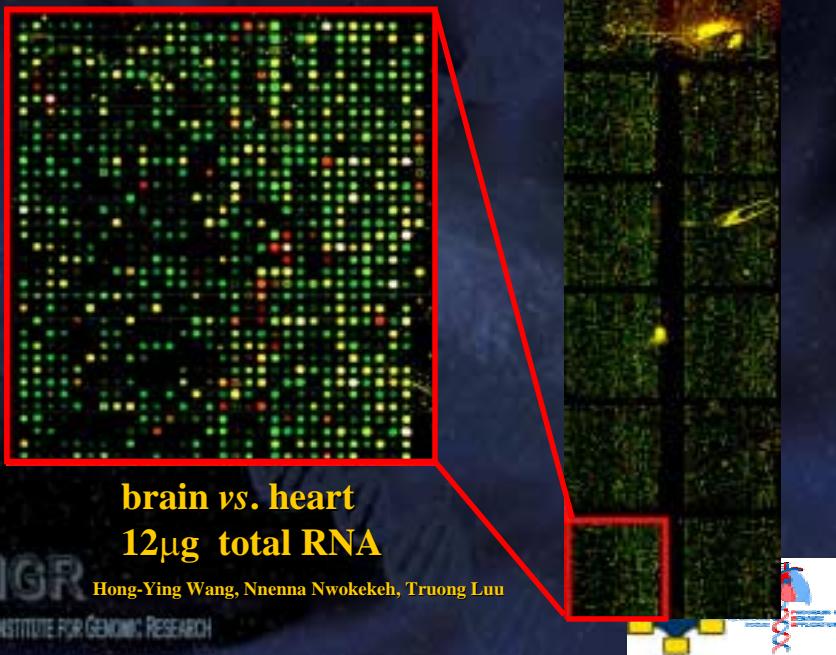
<http://pga.tigr.org/tools.shtml>

## 32,448 element mouse array test



## 13,536 element rat array test



## Spiking Test Samples with *A.thaliana* Control RNAs

Spiked 2-fold change (copies/cell)					Spiked 3-fold change (copies/cell)				
RCA	Cab	rbcL	LTP4	LTP6	XCP2	RPC1	NAC1	TIM	PRK
2 1	10 5	60 30	100 50	300 150	3 1	15 5	60 20	150 50	300 100

**Spike:** Test RNA + Reference RNA → cDNA probe synth. & hybridize

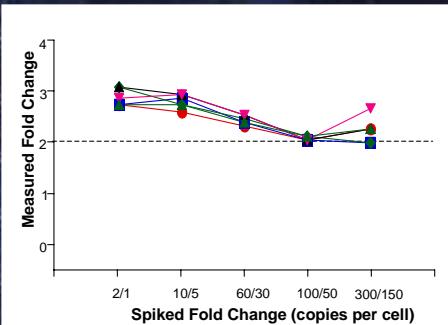
Array containing DNA controls

TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH

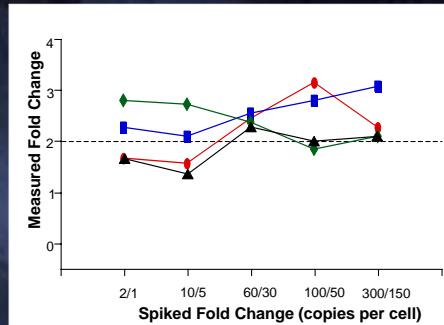


## Assessing Sensitivity of Microarrays at 2:1 ratios

Intra-slide variability across 6 grids:

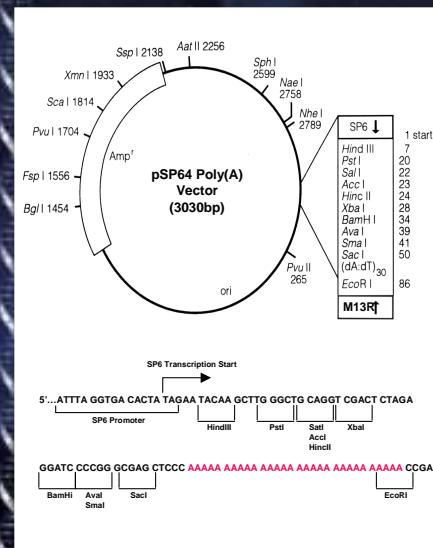


Inter-slide variability:



- Rat cDNA chip
- Total RNA from rat brain and heart
- *A. thaliana* RNA spiking controls: CAB, RCA, rbcL, LTP4, LTP6
- Hong-Ying Wang & Renae Malek

## Resource: *A. thaliana* DNA clones for spiking

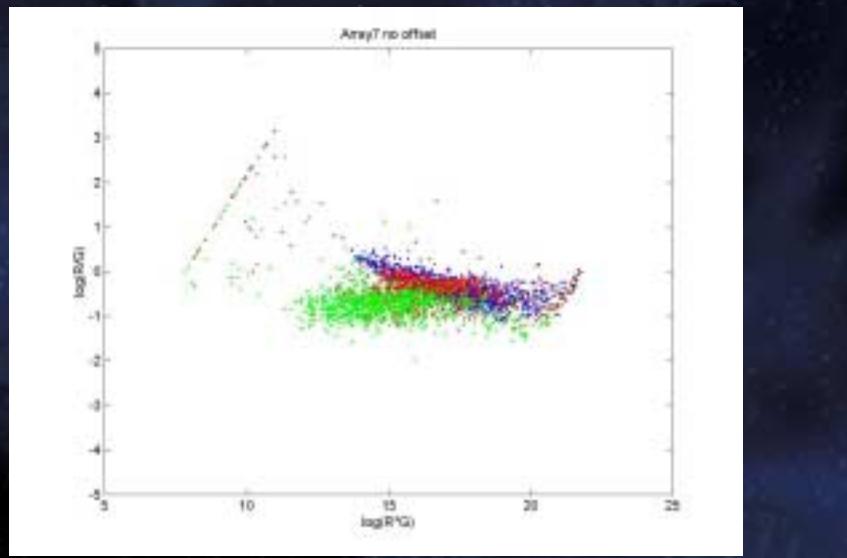


- chlorophyll a/b binding protein (Cab)
- RUBISCO activase (RCA)
- ribulose-1,5-bisphosphate carboxylase/oxygenase (RbcL)
- lipid transfer protein 4 (LTP4)
- lipid transfer protein 6 (LTP6)
- papain-type cysteine endopeptidase (XCP2)
- root cap 1 (RPC1)
- NAC1
- triosphosphate isomerase (TIM)
- ribulose-5-phosphate kinase (PRKase)

TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH



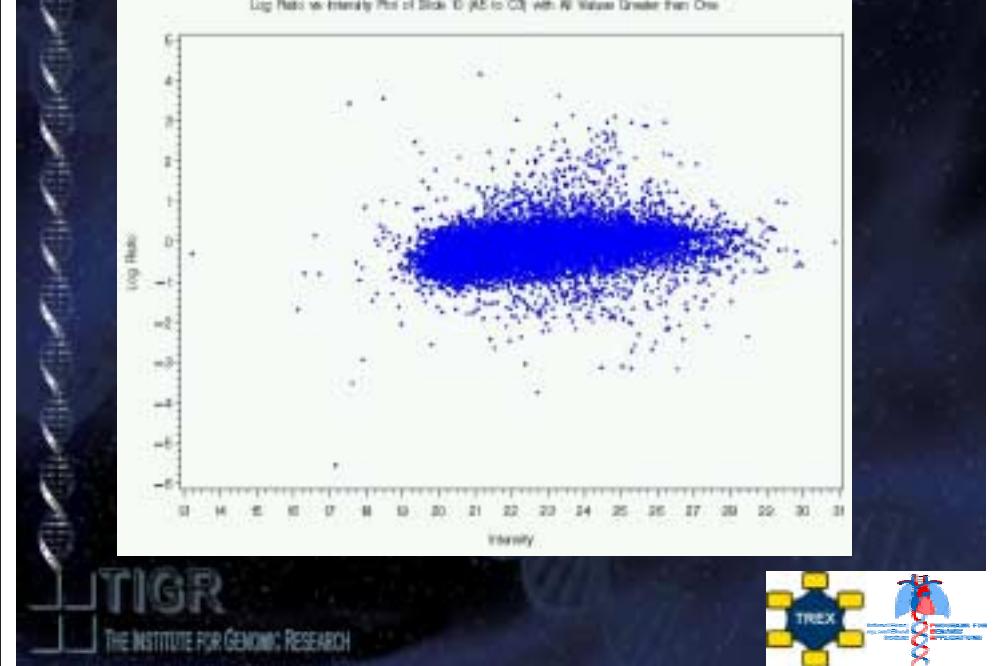
## Bad Data from Parts Unknown



TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH



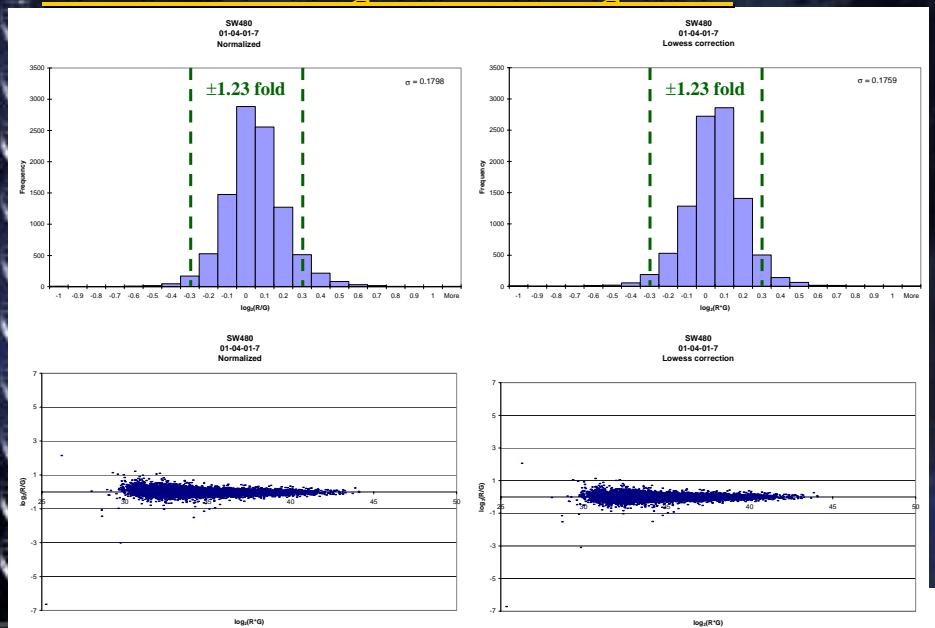
## Good Data from TREX



TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH



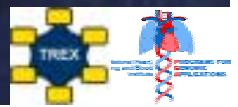
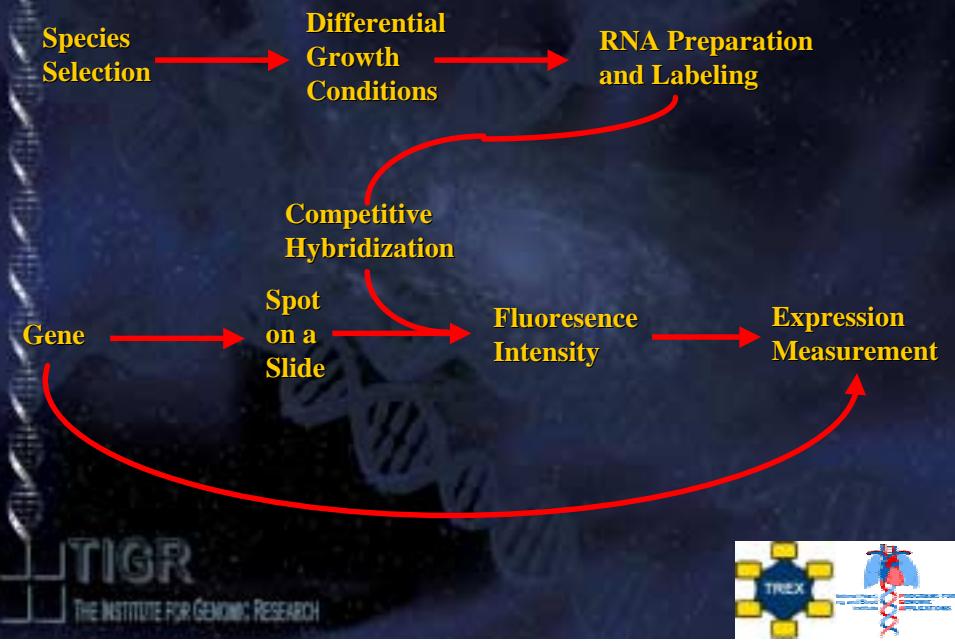
## Normalization using local linear regression



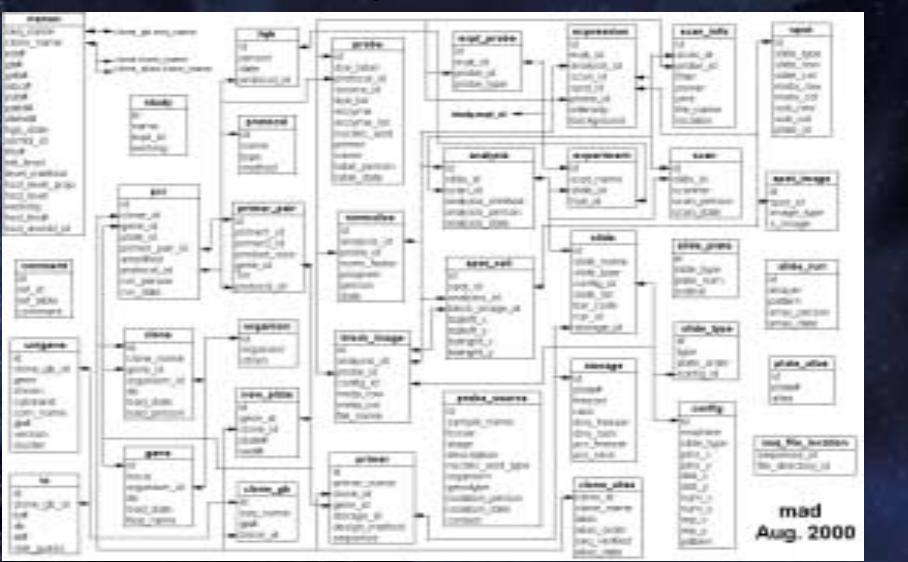
THE INSTITUTE FOR GENOMIC RESEARCH

Ivana Yang , John Quackenbush

## Microarray Expression Analysis



## MAD Microbial Array Database Schema

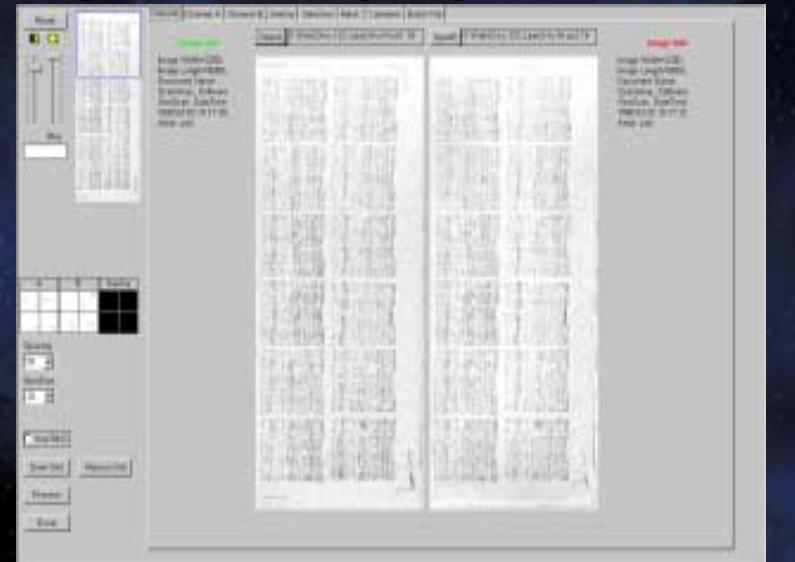


TIGR

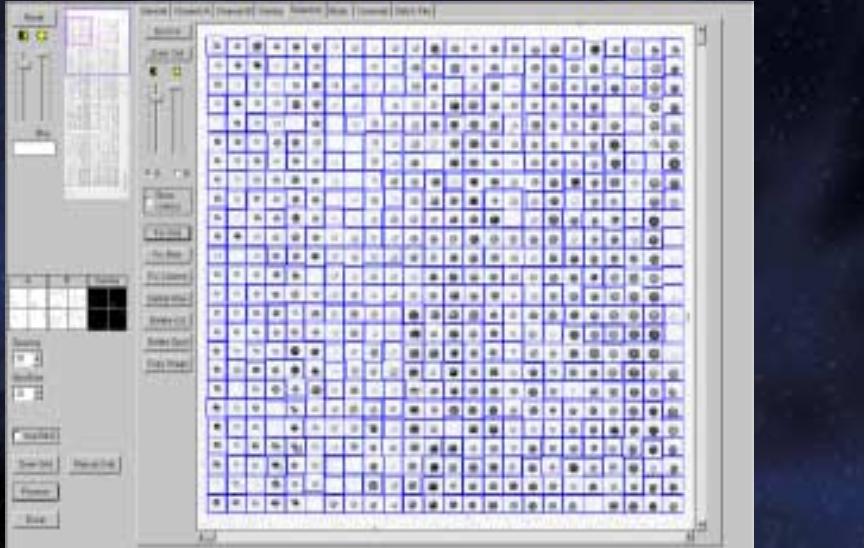
THE INSTITUTE FOR GENOMIC RESEARCH



## TIGR Spotfinder Loading Image Data

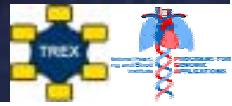


## TIGR Spotfinder Grid Adjustment



TIGR

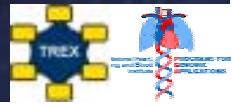
THE INSTITUTE FOR GENOMIC RESEARCH



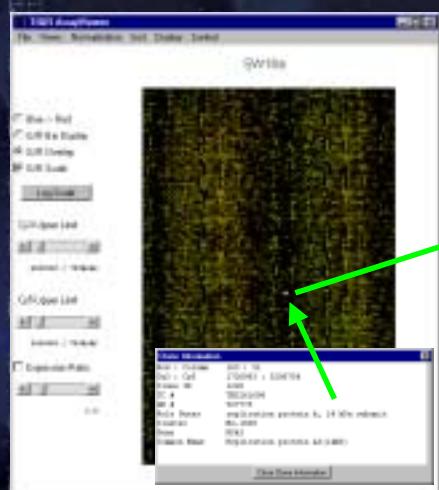
## Data Analysis Issues

- Presentation
- Multiple Views
- Normalization
- Identification of Differentially Expressed Genes
- Multiple Experiments

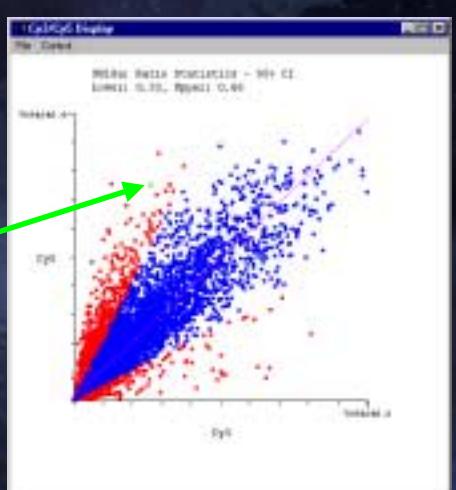
TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH



## TIGR MultiExperiment Viewer

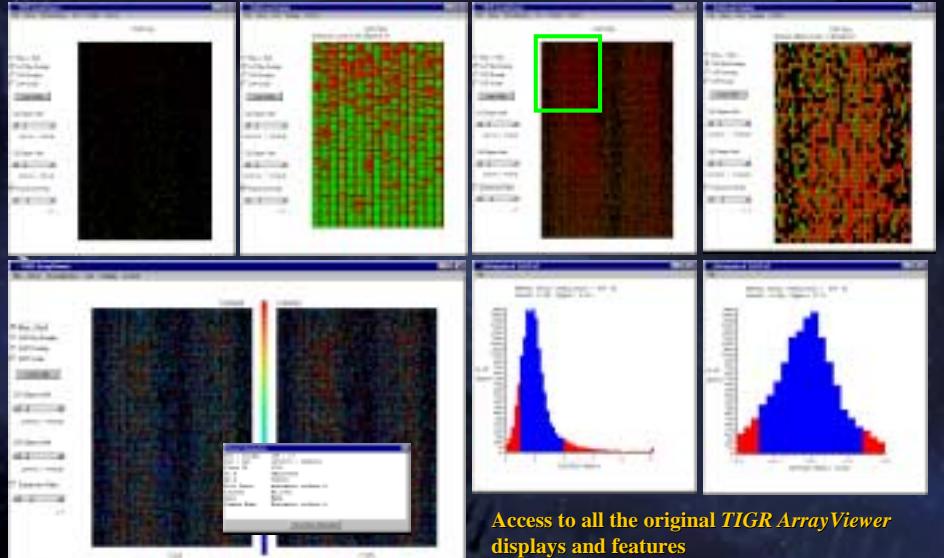


Detail view of a single slide with links to database information about the underlying genes



Highlighted spots in scatter plot track genes being examined in the display on the left

## TIGR MultiExperiment Viewer

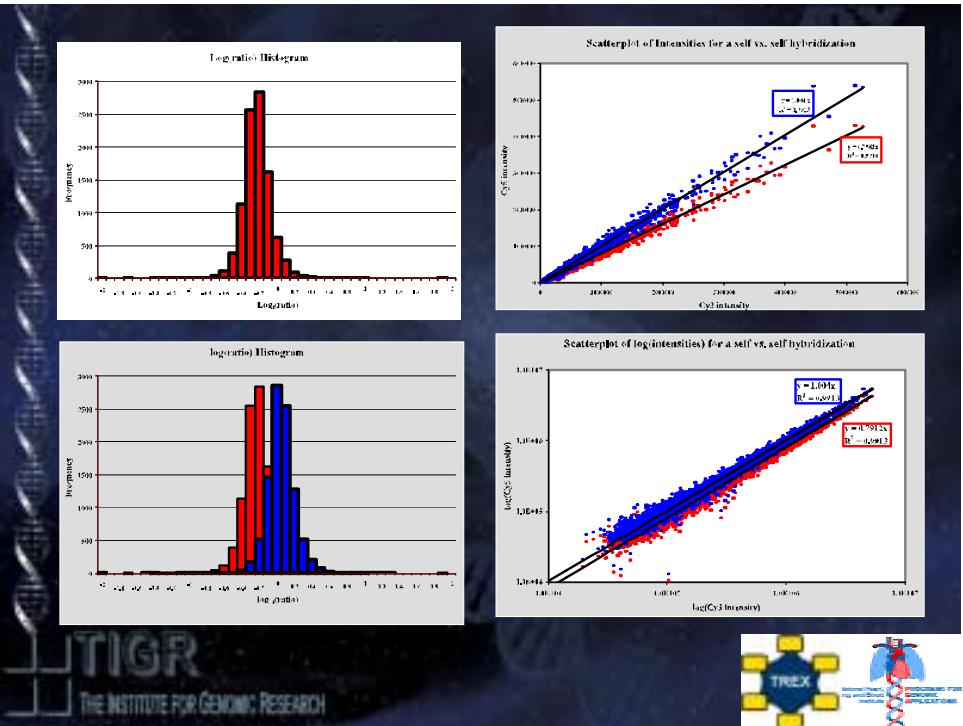


Access to all the original *TIGR ArrayViewer* displays and features



## Why Normalize Data?

- Goal is to measure ratios of gene expression levels  
 $(ratio)_i = R_i/G_i$   
 where  $R_i/G_i$  are, respectively, the measured intensities for the  $i$ th spot.
- In a self-self hybridization, we would expect all ratios to be equal to one:  
 $R_i/G_i = 1$  for all  $i$ . But they may not be.
- Why not?
  - Unequal labeling efficiencies for Cy3/Cy5
  - Noise in the system
  - Differential expression
- Normalization brings (appropriate) ratios back to one.



TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH



## Multiple Experiments?

- Goal is identify genes (or experiments) which have “similar” patterns of expression
- This is a problem in data mining
- “Clustering Algorithms” are most widely used
- Types
  - Agglomerative: Hierarchical
  - Divisive:  $k$ -means, SOMs
  - Others: Principal Component Analysis (PCA)
- All depend on how one measures distance



## Distance metrics

- Distances are measured “between” expression vectors
- Distance metrics define the way we measure distances
- Many different ways to measure distance:
  - Euclidean distance
  - Pearson correlation coefficient(s)
  - Manhattan distance
  - Mutual information
  - Kendall’s Tau
  - etc.
- Each has different properties and can reveal different features of the data

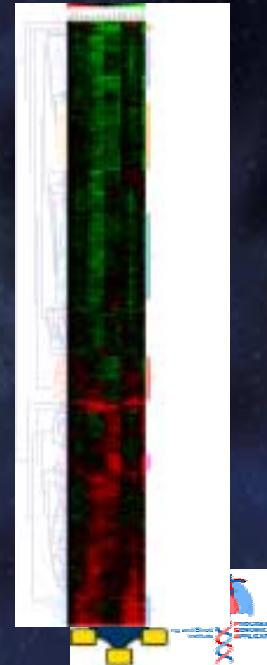
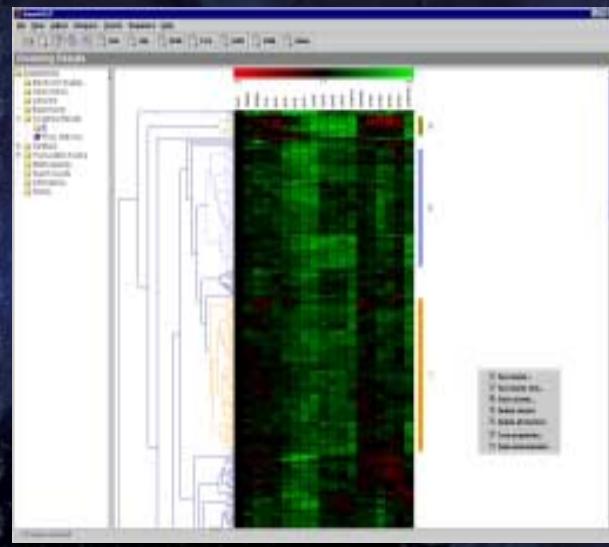


## Expression Vectors

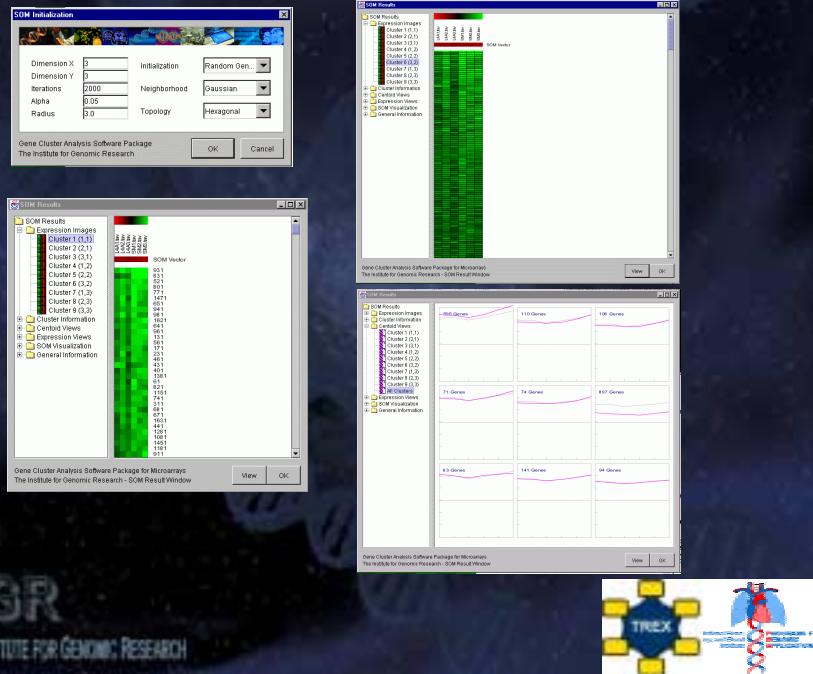
- Crucial concept for understanding clustering
- Each gene is represented by a vector where coordinates are its values  $\log(\text{ratio})$  in each experiment
  - $x = \log(\text{ratio})_{\text{expt1}}$
  - $y = \log(\text{ratio})_{\text{expt2}}$
  - $z = \log(\text{ratio})_{\text{expt3}}$
  - etc.
- For example, if we do six experiments,
  - Gene<sub>1</sub> = (-1.2, -0.5, 0, 0.25, 0.75, 1.4)
  - Gene<sub>2</sub> = (0.2, -0.5, 1.2, -0.25, -1.0, 1.5)
  - Gene<sub>3</sub> = (1.2, 0.5, 0, -0.25, -0.75, -1.4)
  - etc.



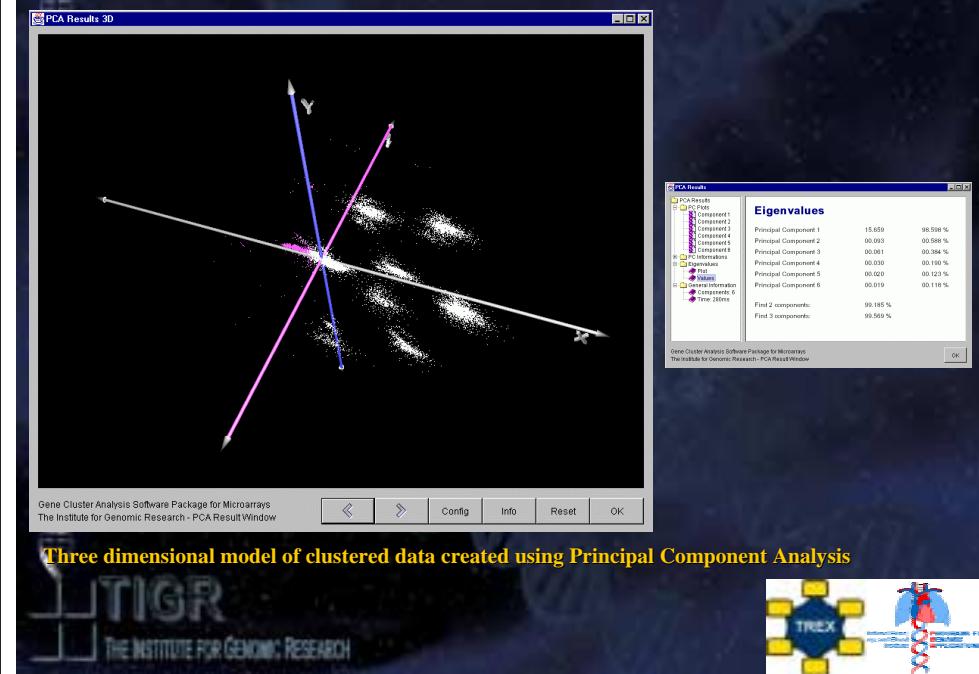
## TIGR MultiExperiment Viewer: Hierarchical Clustering



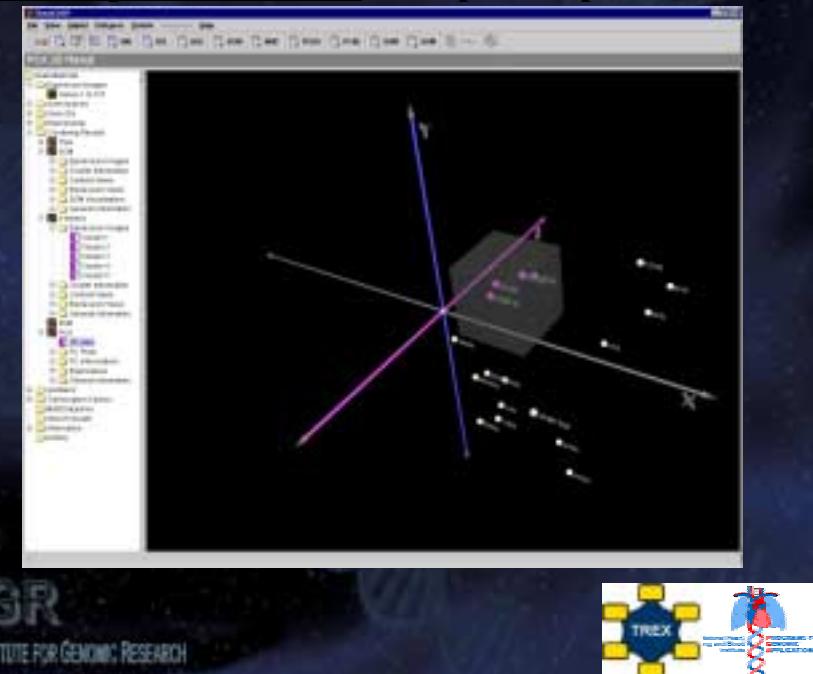
**TIGR MultiExperiment Viewer: Self Organizing Maps**



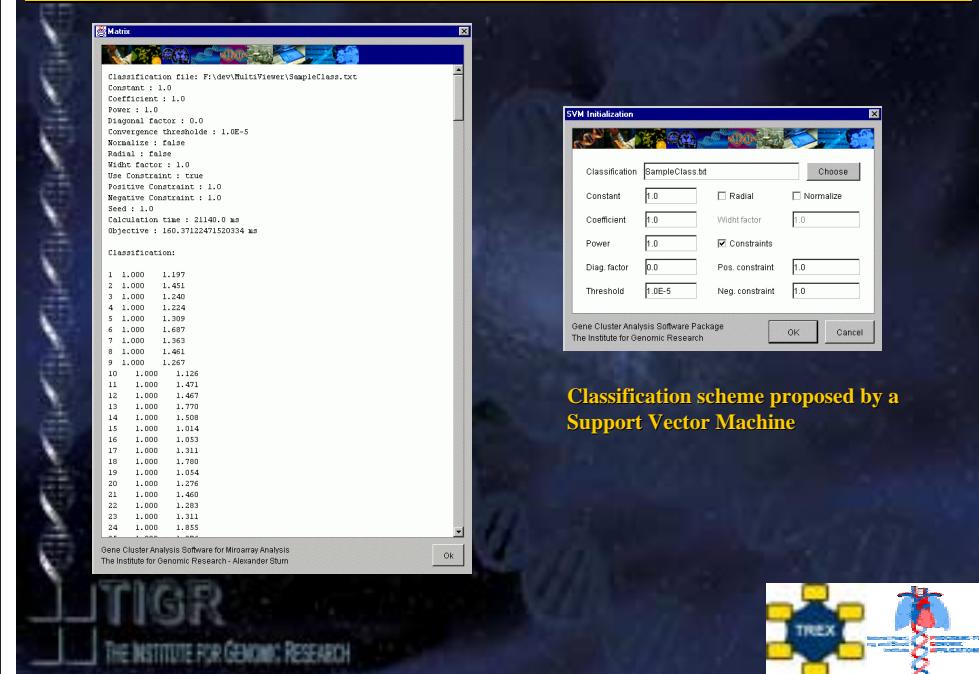
**TIGR MultiExperiment Viewer:Principal Component Analysis**



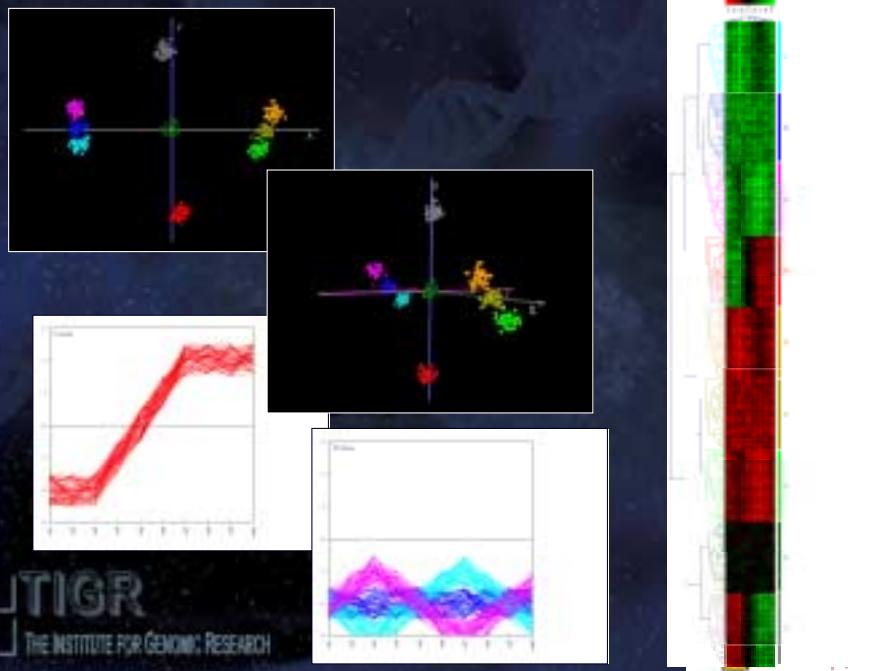
**TIGR MultiExperiment Viewer: Principal Component Analysis**



**TIGR MultiExperiment Viewer: Support Vector Machines**



## TIGR MeV: Tracking across experiments

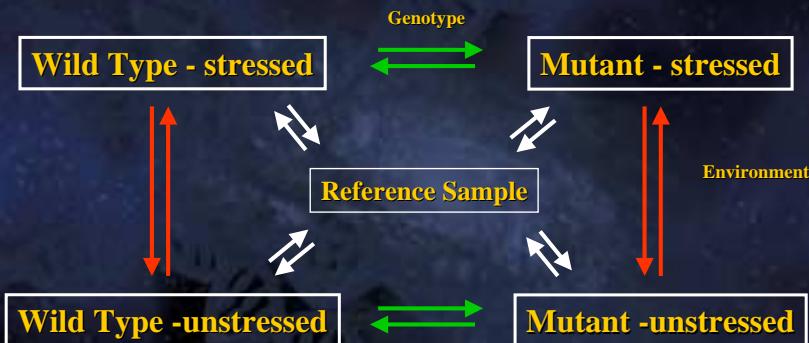


## Developmental Goals

- Tools for linking Genes, ESTs, cDNAs, and Sequences
- Microarray Quality Control Protocols and Reagents
- Novel Analysis Techniques and Tools
- Phenotyping Pipelines
- Microarray Resources in Mouse and Rat
- Expression Profiles for Disease Phenotypes of HLBS interest

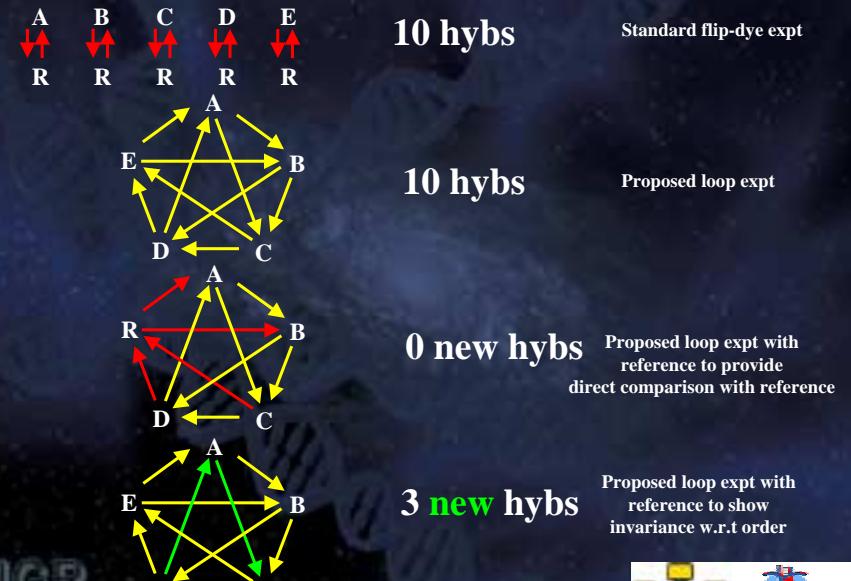


## Basic Experimental Paradigm



TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH

## Loops and Reference Designs

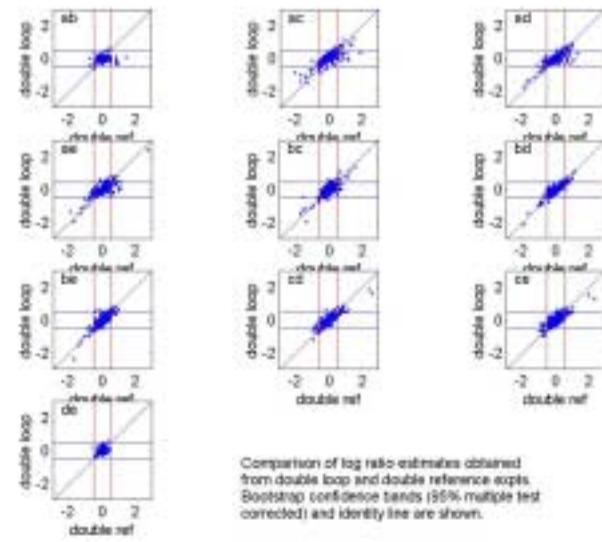


TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH

S. Wang , J. Quackenbush, G. Churchill

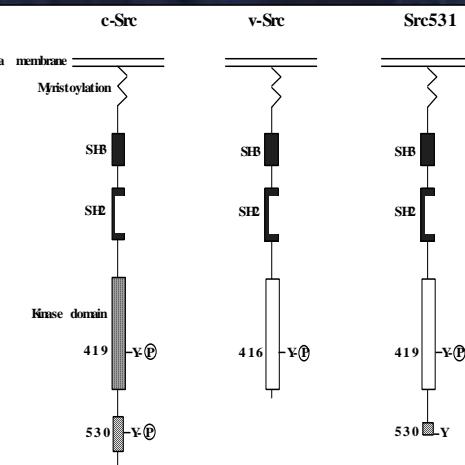


## Loops and Reference Designs



Shuibang Wang, John Quackenbush, Gary Churchill

## An Activating Mutation of Src at Codon 531



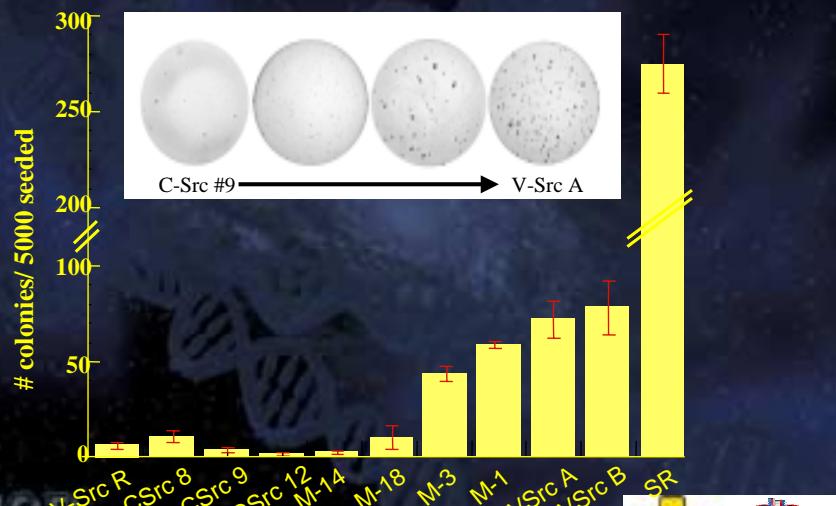
- c-Src  
Non-receptor tyrosine kinase  
Activated by growth factor receptors  
Wide-spread tissue distribution  
Cell growth & differentiation  
Implicated in human colon cancer

- v-Src  
Unregulated kinase activity

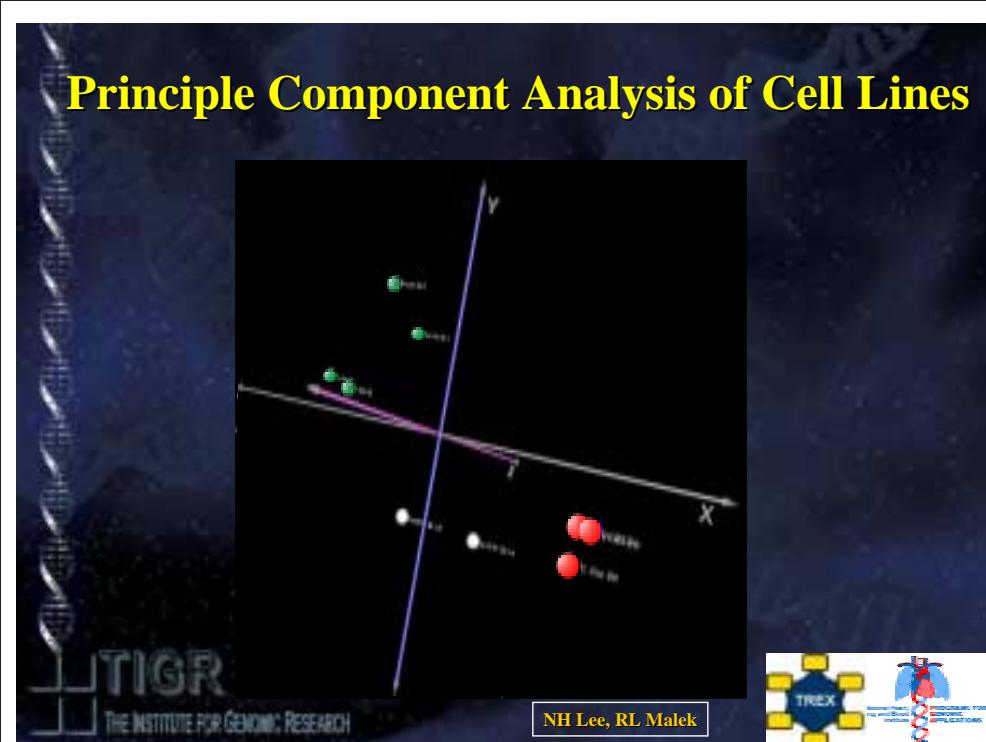
- Src531  
Elevated kinase activity



## Principle Component Analysis of Cell Lines



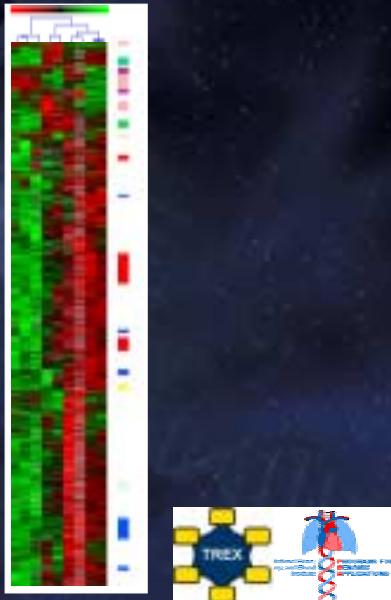
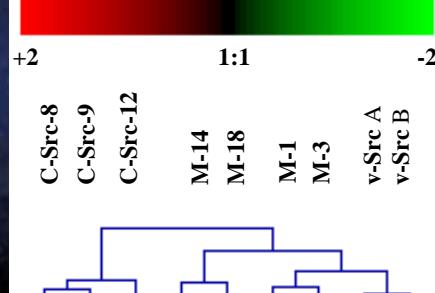
TIGR  
THE INSTITUTE FOR GENOMIC RESEARCH



NH Lee, RL Malek



## Hierarchical Clustering



THE INSTITUTE FOR GENOMIC RESEARCH



## Visiting Scientist Program

**Goal:** To provide access to the reagents, techniques, and analysis tools developed through this PGA

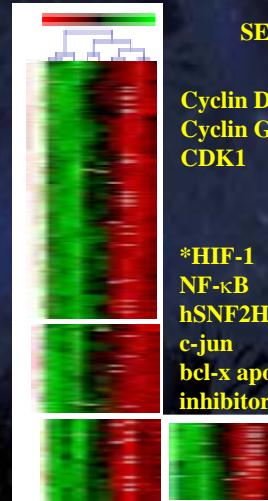
- Solicit applications from external investigators
- PGA steering committee will review applications and select participants
- Seek input from NHLBI staff and PGA-CC
- PGA Investigators will work with Applicants to refine experimental design
- Selected candidates will be invited to spend one (or more) weeks at TIGR generating array data



THE INSTITUTE FOR GENOMIC RESEARCH



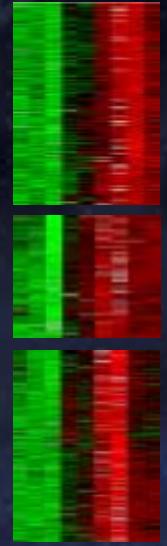
## Transformation Fingerprint



SECIS binding protein 2

Cyclin D1  
Cyclin G  
CDK1  
stanniocalcin  
RYB-a  
arylsulfatase A  
cathepsin B  
hSNF2H  
c-jun  
bcl-x apoptosis  
inhibitor

FA synthase  
stearyl-CoA desaturase  
epoxide hydrolase  
propionyl-CoA carboxylase  
peroxisomal enoyl-CoA hydratase



NH Lee, RL Malek



## Acknowledgments

### The TIGR Gene Index Team

Jennifer Cho  
Svetlana Karamycheva  
Yudan Lee  
Babak Parviz  
Geo Pertea  
Razvan Sultana  
Jennifer Tsai  
John Quackenbush  
Joseph White

*Emeritus*  
Ingeborg Holt (TGI)  
Feng Liang (TGI)  
Kristie Abernathy ( $\mu$ A)  
Sonia Dharap ( $\mu$ A)  
Julie Earle-Hughes ( $\mu$ A)  
Cheryl Gay ( $\mu$ A)  
Priti Hegde ( $\mu$ A)  
Rong Qi ( $\mu$ A)

### H. Lee Moffitt Center/USF

Timothy J. Yeatman

TIGR PGA Collaborators  
Norman Lee  
Renae Malek  
Hong-Ying Wang  
Truong Luu  
Nnenna U. Nwokekeh  
*PGA Collaborators*  
Gary Churchill (TJL)  
Greg Evans (NHLBI)  
Harry Gavaras (BU)  
Howard Jacob (MCW)  
Anne Kwitek-Black (MCW)  
Allan Pack (Penn)  
Beverly Paigen (TJL)  
Luanne Peters (TJL)  
David Schwartz (Duke)

### TIGR Human/Mouse/Arabidopsis Expression Team

Emily Chen  
Renee Gaspard  
Jeremy Hasseman  
Heenam Kim  
John Quackenbush  
Erik Snetsrud  
Shiubang Wang  
Ivana Yang  
Yan Yu  
Baoping Zhao

### Array Software Hit Team

Jerry Li  
John Quackenbush  
Alex Saeed  
Vasily Sharov  
Alexander Sturm  
Joseph White

*Assistant*  
Mary Mulholland

<johnq@tigr.org>

Funding provided by the Department of Energy  
and the National Science Foundation

Funding provided by the National Cancer Institute,  
the National Heart, Lung, Blood Institute,  
and the National Science Foundation

TIGR Faculty, IT Group, and Staff

